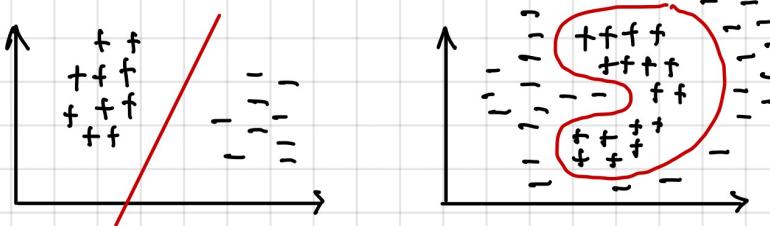


LECTURE 5 : CLASSIFICATION

Motivation: Besides regression, classification is another important factor in making decisions. In fact, most of decisions involve either regression or classification.

Examples of classification :



[1] Classification Linear Model

Data modeling :

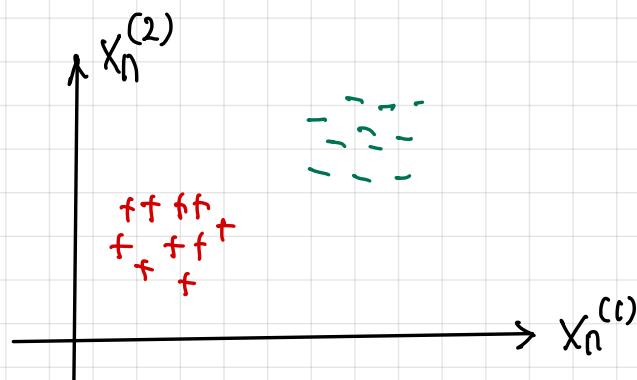
$$y(\bar{x}_n) = \sum_{d=0}^D w_d x_n^{(d)} = w_0 + \sum_{d=1}^D x_n^{(d)}$$

$$= \bar{w}^T \bar{x}_n$$

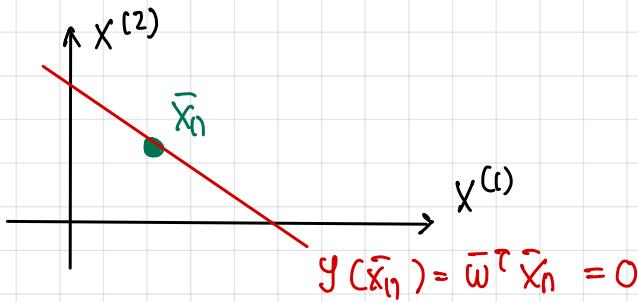
where: $x_n^{(0)} = 1$

$$\bar{x}_n = \begin{bmatrix} 1 \\ x_n^{(1)} \\ x_n^{(2)} \\ \vdots \\ x_n^{(D)} \end{bmatrix}$$

D = Dimensionality of the classification space



Consider a sample \bar{x}_n : $y(\bar{x}_n) = \bar{w}^\top \bar{x}_n = 0$



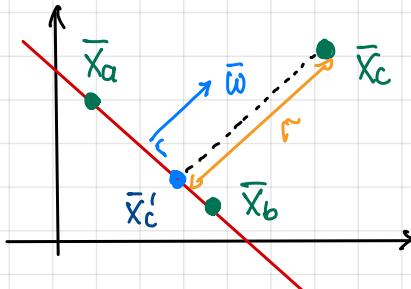
$$\bar{x}_n = \begin{bmatrix} 1 \\ x_n^{(1)} \\ x_n^{(2)} \end{bmatrix}$$

$$y(\bar{x}_n) = \bar{w}^\top \bar{x}_n = [w_0 \ w_1 \ w_2] \begin{bmatrix} 1 \\ x_n^{(1)} \\ x_n^{(2)} \end{bmatrix} = 0$$

$$= w_0 + w_1 x_n^{(1)} + x_n^{(2)} = 0$$

Recall: $ax + by + c = 0 \rightarrow y = -\frac{a}{b}x - \frac{c}{b}$: line equation

Consider two samples \bar{x}_a and \bar{x}_b :



$$y(\bar{x}_a) = y(\bar{x}_b) = 0$$

$$y(\bar{x}_a) - y(\bar{x}_b) = 0$$

$$\bar{w}^\top (\bar{x}_a - \bar{x}_b) = 0$$

Recall: $\bar{a}^\top \bar{b} = \bar{a} \cdot \bar{b} = \|\bar{a}\| \|\bar{b}\| \cos \theta$
if $\bar{a}^\top \bar{b} = 0$, then $\theta = 90^\circ$

$(\bar{x}_a - \bar{x}_b)$ is a line segment lying on the red line, $y(\bar{x}_a) = 0$.

Hence: \bar{w} is perpendicular to the red line.

Assume \bar{x}_c , which is not on the line: $y(\bar{x}_c) \neq 0$

$$\bar{x}_c = \bar{x}_c' + r \frac{\bar{w}}{\|\bar{w}\|}$$

$$y(\bar{x}_c) = \bar{w}^\top \left(\bar{x}_c' + r \frac{\bar{w}}{\|\bar{w}\|} \right)$$

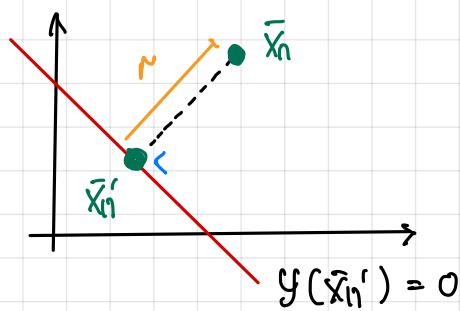
$$= \bar{w}^\top \bar{x}_c' + \bar{w}^\top r \frac{\bar{w}}{\|\bar{w}\|}$$

$$= 0 + r \frac{\bar{w}^\top \bar{w}}{\|\bar{w}\|}$$

$$\|\bar{w}\| \|\bar{w}\| \cos 0 \\ = \|\bar{w}\|^2$$

$$y(\bar{x}_c) = r \frac{\bar{w}^\top \bar{w}}{\|\bar{w}\|} = r \frac{\|\bar{w}\|^2}{\|\bar{w}\|} = r \|\bar{w}\|$$

If $y(\bar{x}_n) = 1$, then: $y(\bar{x}_n) = \bar{w}^T \bar{x}_n = r \|\bar{w}\| = 1$



$$r = 1/\|\bar{w}\|$$

↓

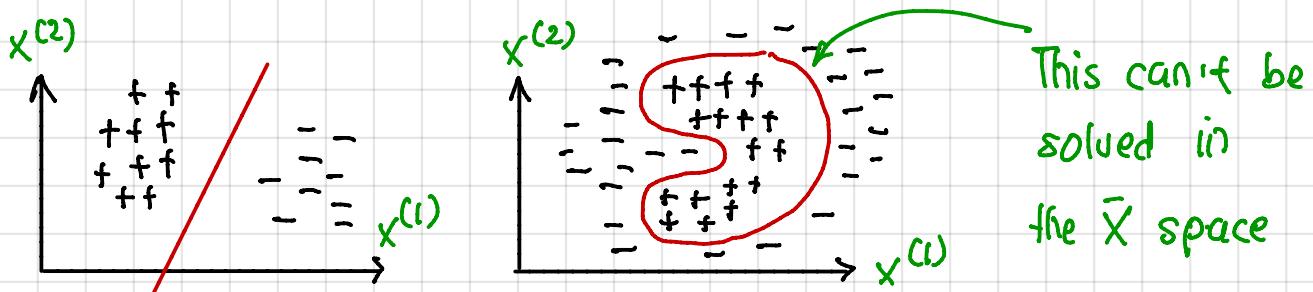
$y(\bar{x}_n) = 1$ implies \bar{x}_n is located outside the line with distance r .

SPACE TRANSFORMATION

So far, we define: $\bar{Y}(\bar{x}_n) = \bar{W}^T \bar{\phi}(\bar{x}_n)$

This comes with a drawback:

Our boundary line is always linear, however the distributions of the data belonging to class 1 & 2 might not be linear:



To address this problem, we need to transform \bar{x}_n into another (usually higher) space:

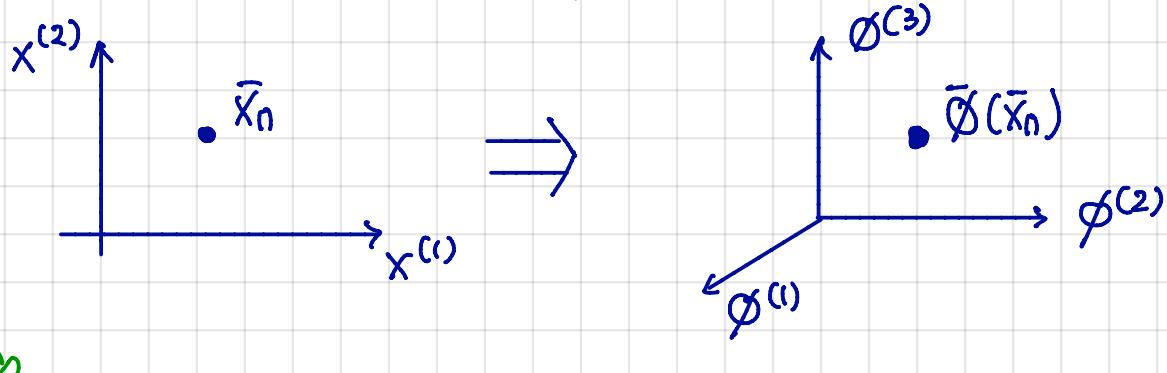
$$y_k(\bar{x}_n) = \bar{W}_k^T \bar{\phi}(\bar{x}_n) = \sum_{m=0}^M w_k^{(m)} \phi^{(m)}(\bar{x}_n)$$

$|x|_1 \quad |x|_M \quad M|x|_1$

e.g.: $\phi^{(m)}(\bar{x}_n) = \exp\left(-\frac{\|\bar{x}_n - \bar{\mu}_m\|_2^2}{2s^2}\right)$; $\|\bar{a}\|_2^2 = \sqrt{a_1^2 + a_2^2}$

$\bar{\mu}_m$ is a vector with the same elements/values:

$$\|\bar{x}_n - \bar{\mu}_m\|_2^2 = (x_n^{(1)} - \mu_m)^2 + (x_n^{(2)} - \mu_m)^2 \quad ; \quad \bar{\mu}_m \text{ is the centers of clusters}$$



$$\bar{Y}(\bar{x}_n) = \bar{W}^T \bar{\phi}(\bar{x}_n) \quad \rightarrow \quad \bar{Y} = \bar{\phi} \bar{W}$$

$|x|_1 \quad |x|_M \quad M|x|_1$

$$N|x|_K \quad N|X|M \quad M|x|_K$$

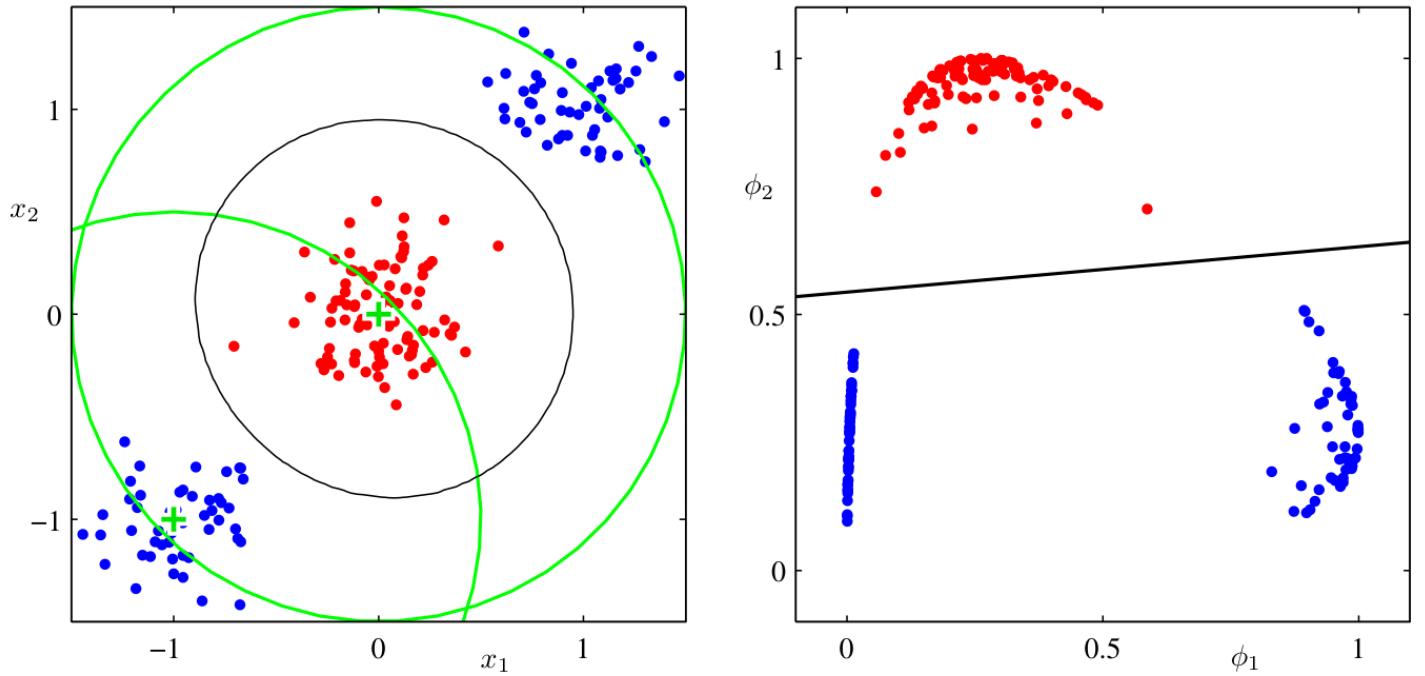


Figure 4.12 Illustration of the role of nonlinear basis functions in linear classification models. The left plot shows the original input space (x_1, x_2) together with data points from two classes labelled red and blue. Two ‘Gaussian’ basis functions $\phi_1(\mathbf{x})$ and $\phi_2(\mathbf{x})$ are defined in this space with centres shown by the green crosses and with contours shown by the green circles. The right-hand plot shows the corresponding feature space (ϕ_1, ϕ_2) together with the linear decision boundary obtained given by a logistic regression model of the form discussed in Section 4.3.2. This corresponds to a nonlinear decision boundary in the original input space, shown by the black curve in the left-hand plot.

MLE FOR CLASSIFICATION

(LOGISTIC REGRESSION)

Testing stage: Input: \bar{x}_* and \bar{w}

Output: t_* (label of \bar{x}_*) ; $t_* \in \{0, 1\}$

Goal: $p(t_* | \bar{x}_*, \bar{w})$



Assuming two classes only

$0 = \text{class 1}$

$1 = \text{class 2}$

(1) We want to know $p(t_* = 0 | \bar{x}_*, \bar{w})$ and $p(t_* = 1 | \bar{x}_*, \bar{w})$

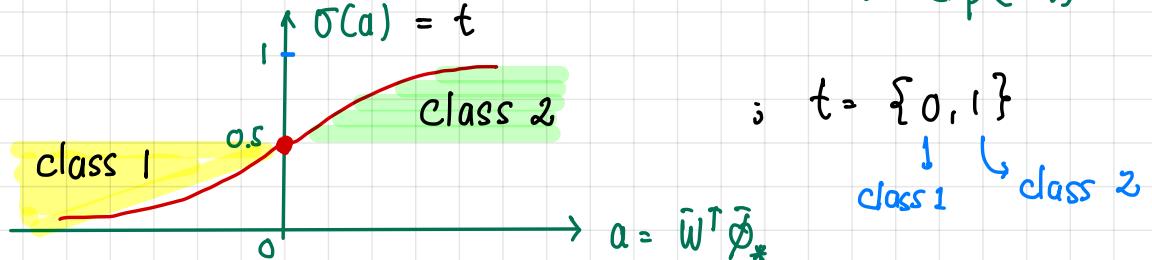
(2) if $p(t_* = 0 | \bar{x}_*, \bar{w}) \geq p(t_* = 1 | \bar{x}_*, \bar{w})$

then $t_* = 0 \rightarrow \bar{x}_*$ belongs to class 1

else $t_* = 1 \rightarrow \bar{x}_*$ belongs to class 2

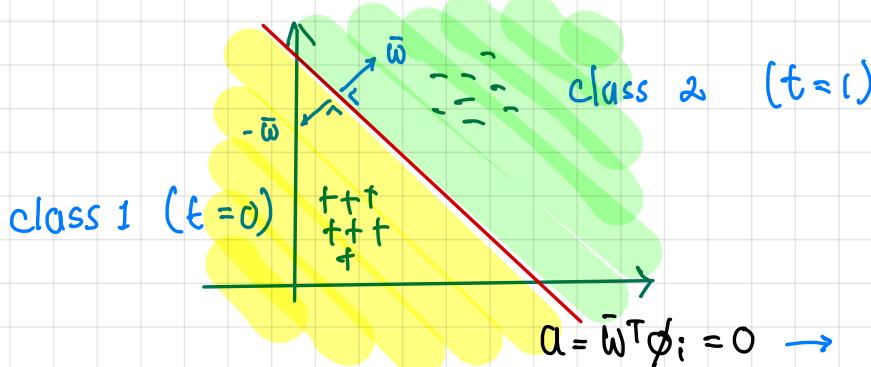
Q: $p(t_* | \bar{x}_*, \bar{w}) = ?$ How can we model the probability function?

A: $p(t_* | \bar{x}_*, \bar{w}) = \sigma(\bar{w}^T \phi(\bar{x}_*))$; $\sigma(a) = \frac{1}{1 + \exp(-a)}$



$t_* = 0$: $a = \bar{w}^T \phi_*$ is encouraged to be negative & large : $-a = -\bar{w}^T \phi_*$

$t_* = 1$: $a = \bar{w}^T \phi_*$ is encouraged to be positive & large : $a = \bar{w}^T \phi_*$



we only have one set of \bar{w}

Training stage: Input: $\{\bar{x}_n, t_n\}_{n=1}^N$

Output: \bar{w}

(1) Likelihood data modeling:

Likelihood: $p(t_n | \bar{w})$

→ One sample

$$p(t_n = 1 | \bar{w}) = \sigma(\bar{w}^T \bar{\phi}_n) \\ = y_n$$

$$p(t_n = 0 | \bar{w}) = 1 - \sigma(\bar{w}^T \bar{\phi}_n) \\ = 1 - y_n$$

To make them into 1 equation:

$$p(t_n | \bar{w}) = y_n^{t_n} (1-y_n)^{1-t_n}$$

For all N data:

$$p(\bar{t} | \bar{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$$

Q: How do we know that it's a probability function?

A: For simplicity, let N=2:

$$p(t_1=1, t_2=1 | \bar{w}) = y_1 y_2$$

$$p(t_1=1, t_2=0 | \bar{w}) = y_1 (1-y_2) = y_1 - y_1 y_2$$

$$p(t_1=0, t_2=1 | \bar{w}) = (1-y_1) y_2 = y_2 - y_1 y_2$$

$$p(t_1=0, t_2=0 | \bar{w}) = (1-y_1) (1-y_2) = 1 - y_1 - y_2 + y_1 y_2 +$$

$$\sum_{t_1} \sum_{t_2} p(t_1, t_2 | \bar{w}) = 1$$

(2) Error function:

$$\bar{\omega}^* = \arg \max_{\{\bar{\omega}\}} p(\bar{t} | \bar{\omega}) = \arg \min_{\{\bar{\omega}\}} E(\bar{\omega})$$

$$\begin{aligned} E(\bar{\omega}) &= -\log p(\bar{t} | \bar{\omega}) \\ &= -\sum_n (t_n \log y_n + (1-t_n) \log (1-y_n)) \\ &= -\sum_n [t_n \log \sigma(\bar{\omega}^\top \bar{\phi}_n) + (1-t_n) \log (1-\sigma(\bar{\omega}^\top \bar{\phi}_n))] \end{aligned}$$

(3) Minimization:

$$\frac{\partial E(\bar{\omega})}{\partial \bar{\omega}} = 0 = -\sum_n \frac{\partial}{\partial \bar{\omega}} [t_n \log \sigma + (1-t_n) \log (1-\sigma)]$$

Note: $\frac{\partial \sigma(x)}{\partial x} = \sigma(1-\sigma)$

$$\frac{\partial (1-\sigma(x))}{\partial x} = -\sigma'(1-\sigma)$$

$\left. \begin{array}{l} \text{see the last} \\ \text{page for the} \\ \text{proof.} \end{array} \right\}$

$$\begin{aligned} &= -\sum_n t_n \cancel{\frac{\sigma(1-\sigma)}{\sigma}} \bar{\phi}_n - \sum_n (1-t_n) \cancel{\frac{\sigma(1-\sigma)}{1-\sigma}} \bar{\phi}_n = 0 \\ &= \sum_n -t_n \bar{\phi}_n + \cancel{t_n \sigma \bar{\phi}_n} + \sigma \bar{\phi}_n - \cancel{t_n \sigma \bar{\phi}_n} = 0 \end{aligned}$$

$$\frac{\partial E(\bar{\omega})}{\partial \bar{\omega}} = \boxed{\nabla E(\bar{\omega}) = \sum_n (\sigma - t_n) \bar{\phi}_n} = 0$$

Q: How can we get $\bar{\omega}$ from the last equation?

A: Unfortunately, we cannot solve it in a closed-form manner.
Because $\sigma(-\bar{\omega}^\top \bar{\phi}_n)$ is not linear.

[•] Iterative Reweighted Least Squares

Newton's method (or Newton-Raphson method):

$$\bar{w}^{\text{new}} = \bar{w}^{\text{old}} - H^{-1} \nabla E(\bar{w})$$

where: H is the Hessian matrix, which is $\frac{\partial^2 E(\bar{w})}{\partial \bar{w}^2}$ or $\frac{\partial \nabla E(\bar{w})}{\partial \bar{w}}$

$$\begin{aligned} \nabla E(\bar{w}) &= \sum_n (y_n - t_n) \phi_n = \sum_n (\sigma(\bar{w}^T \phi_n) - t_n) \phi_n \\ &= \underbrace{\Phi^T}_{M \times N} \underbrace{(y - t)}_{N \times 1} \end{aligned}$$

$$\begin{aligned} H &= \frac{\partial}{\partial \bar{w}} \left(\sum_n (\sigma_n - t_n) \bar{\phi}_n \right) = \sum_n \underbrace{\bar{\phi}_n}_{M \times 1} \underbrace{\sigma_n}_{1 \times 1} \underbrace{(1 - \sigma_n)}_{1 \times 1} \underbrace{\bar{\phi}_n^T}_{1 \times M} \\ &= \underbrace{\Phi^T}_{M \times N} \underbrace{R}_{N \times N} \underbrace{\Phi}_{N \times M} \end{aligned}$$

where : $R_{nn} = \sigma_n(1 - \sigma_n)$

R is a diagonal matrix



$$\bar{w}^{\text{new}} = \bar{w}^{\text{old}} - H^{-1} \nabla E(\bar{w})$$

$$= \bar{w}^{\text{old}} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

$$= (\Phi^T R \Phi)^{-1} [\Phi^T R \Phi \bar{w}^{\text{old}} - \Phi^T (y - t)]$$

$$= (\Phi^T R \Phi)^{-1} \Phi^T R \bar{z} \quad \rightarrow \text{least-squares solution:}$$

$$(\Phi^T \Phi)^{-1} \Phi^T \bar{t}$$

where : $\bar{z} = \Phi \bar{w}^{\text{old}} - R^{-1} (y - t)$

Hence the name :

Iterative reweighted least-squares

For our binary classification case we can initially set $\bar{w} = 0$

Note :

$$\begin{aligned}
 \frac{d\sigma(x)}{dx} &= \frac{d}{dx} \left(\frac{1}{1 + \exp(-x)} \right) = \frac{d}{dx} \left(1 + \exp(-x) \right)^{-1} \\
 &= - \left(1 + \exp(-x) \right)^{-2} (-\exp(-x)) = \frac{\exp(-x)}{(1 + \exp(-x))^2} = \sigma^2 \exp(-x) \\
 &= \sigma^2 \left(1 + \exp(-x) - 1 \right) = \sigma \left(\frac{1 + \exp(-x) - 1}{1 + \exp(-x)} \right) \\
 &= \sigma \left(1 - \frac{1}{1 + \exp(-x)} \right) \\
 &= \sigma (1 - \sigma)
 \end{aligned}$$

Accordingly : $\frac{d}{dx} ((1-\sigma)) = -\sigma (1-\sigma) = \sigma (\sigma-1)$

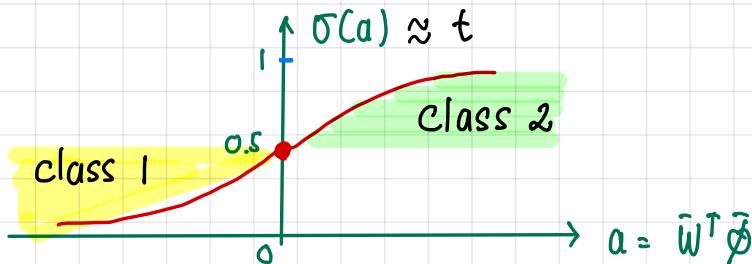
MAP FOR CLASSIFICATION

Previously we know : $t_n \in \{0, 1\}$

$$t_n \approx y_n = \sigma(\bar{w}^\top \bar{\phi}(x_n)) = \sigma(\bar{w}^\top \bar{\phi}_n)$$

(XM) (MX1)

(XM) (MX1)



if $t_n=0$: $a = \bar{w}^\top \bar{\phi}_n$ is encouraged to be large & negative

if $t_n=1$: $a = \bar{w}^\top \bar{\phi}_n$ is encouraged to be large & positive

Unlike least squares, we have only 1 line, as we have only 1 set of \bar{w} .

- Training Stage using MAP :

Input: $\{\bar{x}_n, t_n\}_{n=1}^N$
Output: \bar{w}

(I) Posterior Data Modeling :

$$\text{posterior: } p(\bar{w} | \bar{t}) \propto p(\bar{t} | \bar{w}) p(\bar{w})$$

(MX1) (NX1) (MX1) (MX1)

likelihood prior

$$\text{where: likelihood } p(\bar{t} | \bar{w}) = \prod_n p(t_n | \bar{w})$$

$$= \prod_n y_n^{t_n} (1-y_n)^{1-t_n}$$

$$\text{Prior } p(\bar{w}) = G(\bar{w}; \bar{M}_0, \bar{S}_0)$$

(MX1) (MX1) (MXM)

$\bar{M}_0 = 0$
 $\bar{S}^{-1} = \alpha^{-1} \mathbb{I}$

see last page for more discussion

Hence :

$$\text{posterior } p(\bar{w} | \bar{t}) \propto p(\bar{t} | \bar{w}) p(\bar{w}) = \prod_n y_n^{t_n} (1-y_n)^{1-t_n} G(\bar{w}; \bar{M}_0, \bar{S}_0)$$

(2) MAP Error Function:

$$\hat{\vec{w}}_{\text{MAP}} = \underset{\{\vec{w}\}}{\operatorname{argmax}} p(\vec{w} | \vec{t}) = \underset{\{\vec{w}\}}{\operatorname{argmin}} E(\vec{w})$$

where:

$$E(\vec{w}) = -\log p(\vec{w} | \vec{t}) ; \quad \vec{o}_n = \sigma(\vec{w}^T \vec{\phi}_n) = y_n$$

$$= - \left[\sum_n t_n \log \sigma_n + (1-t_n) \log (1-\sigma_n) \right] + \frac{1}{2} (\vec{w} - \vec{m}_0)^T \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) + \text{const}$$

$$(3) \text{ Minimization: } \frac{\partial E(\vec{w})}{\partial \vec{w}} = 0$$

MX1

$$\begin{aligned} \frac{\partial E(\vec{w})}{\partial \vec{w}} &= - \left[\sum_n t_n \cancel{\frac{1}{\sigma_n}} (-\sigma_n) \vec{\phi}_n - \sum_n (1-t_n) \sigma_n \cancel{\frac{-1}{(1-\sigma_n)}} \vec{\phi}_n \right] + \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) \\ &\stackrel{\text{MX1}}{=} - \sum_n t_n (1-\sigma_n) \vec{\phi}_n + \sum_n (1-t_n) \sigma_n \vec{\phi}_n + \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) \\ &= - \sum_n t_n \vec{\phi}_n + \sum_n t_n \sigma_n \vec{\phi}_n + \sum_n \sigma_n \vec{\phi}_n - \sum_n t_n \sigma_n \vec{\phi}_n + \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) \\ &= - \left[\sum_n (t_n \vec{\phi}_n - \sigma_n \vec{\phi}_n) \right] + \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) \\ &= \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) + \sum_n (\sigma_n - t_n) \vec{\phi}_n \end{aligned}$$

$$\boxed{\nabla E(\vec{w}) = \frac{\partial E(\vec{w})}{\partial \vec{w}} = \mathbb{S}_0^{-1} (\vec{w} - \vec{m}_0) + \vec{\phi}^T \frac{(\vec{o} - \vec{t})}{\text{NEX1}}}$$

$$\boxed{H = \nabla \nabla E(\vec{w}) = \mathbb{S}_0^{-1} + \sum_n \sigma_n (1-\sigma_n) \vec{\phi}_n \vec{\phi}_n^T}$$

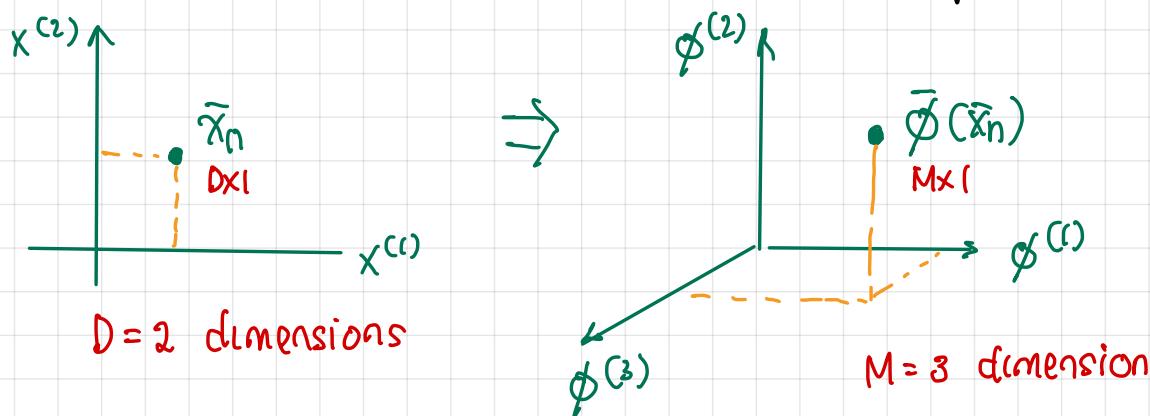
$$H = \mathbb{S}_0^{-1} + \vec{\phi} R \vec{\phi}^T \quad . \quad R_{nn} = \sigma_n (1-\sigma_n)$$

$$\text{Newton's method: } \vec{w}^{\text{new}} = \vec{w}^{\text{old}} - H^{-1} \nabla E(\vec{w})$$

Notes :

1. Q: Why do we model $p(\bar{w}) = \mathcal{G}(\bar{w}; \theta, \propto \mathbb{I})$? What is the meaning of this prior?

A: In regression, the model complexity M determines the complexity of the fitting line. In classification, M determines the dimensionality of the transform space. The classification line / plane is always linear.



Hence, large M implies a space with high dimensionality, where we can always separate the positive from negative classes. This implies overfitting.



To reduce the overfitting problem, we encourage some of w's to be zero (= reducing the dimensions)

FULL BAYESIAN CLASSIFICATION

Training stage : Input: $\{\bar{x}_n, t_n\}_{n=1}^N$; $\bar{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}$
 Output: \bar{w}

(i) Full Bayesian Data Modeling : Likelihood Prior

$$\text{Posteriori : } p(\bar{w}|\bar{t}) = \frac{p(\bar{t}|\bar{w}) p(\bar{w})}{p(\bar{t})}$$

$$\text{Likelihood : } p(\bar{t}|\bar{w}) = \prod_n \sigma_n^{t_n} (r - \sigma_n)^{1-t_n}$$

$$\text{Prior : } p(\bar{w}) = G(\bar{w}; \bar{m}_0, S_0)$$

Hence :

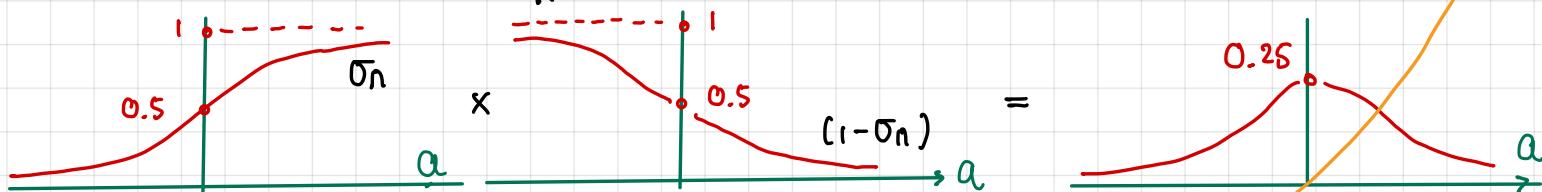
$$p(\bar{w}|\bar{t}) = \frac{\prod_n \sigma_n^{t_n} (r - \sigma_n)^{1-t_n} G(\bar{w}; \bar{m}_0, S_0)}{\int \prod_n \sigma_n^{t_n} (r - \sigma_n)^{1-t_n} G(\bar{w}; \bar{m}_0, S_0) d\bar{w}}$$

$$\int \prod_n \sigma_n^{t_n} (r - \sigma_n)^{1-t_n} G(\bar{w}; \bar{m}_0, S_0) d\bar{w}$$

Problems :

1. In regression, we can solve the full Bayesian equation using Case 2.
 Case 2 requires both likelihood & prior to be Gaussians. Yet, in the above equation, the likelihood is NOT Gaussian.

2. Likelihood : $p(\bar{t}|\bar{w}) = \prod_n \sigma_n^{t_n} (r - \sigma_n)^{1-t_n}$ is not Gaussian :



$$3. f(\bar{w}) = \underbrace{\prod_n \sigma_n^{t_n} (r - \sigma_n)^{1-t_n}}_{\text{looks like a Gaussian}} G(\bar{w}; \bar{m}_0, S_0)$$

looks like
a Gaussian

[•] Laplace Approximation

Goal : To approximate $f(\bar{w})$ with a Gaussian function.

$$f(\bar{w}) = \prod_n \sigma_n t^n (r - \sigma_n)^{c_r - t_n} G(\bar{w}; \bar{m}_0, \mathbb{S}_0)$$

$$\approx f(\bar{w}) = G(\bar{w}; \bar{m}_N, \mathbb{S}_N)$$

to find

Concept :

$$p(x) = \frac{1}{z} f(x)$$

$$; z = \int f(x) dx$$

$$q(x) \approx p(x)$$

unknown normalization factor

To find : $q(x)$, a Gaussian approximation

Step 1 :

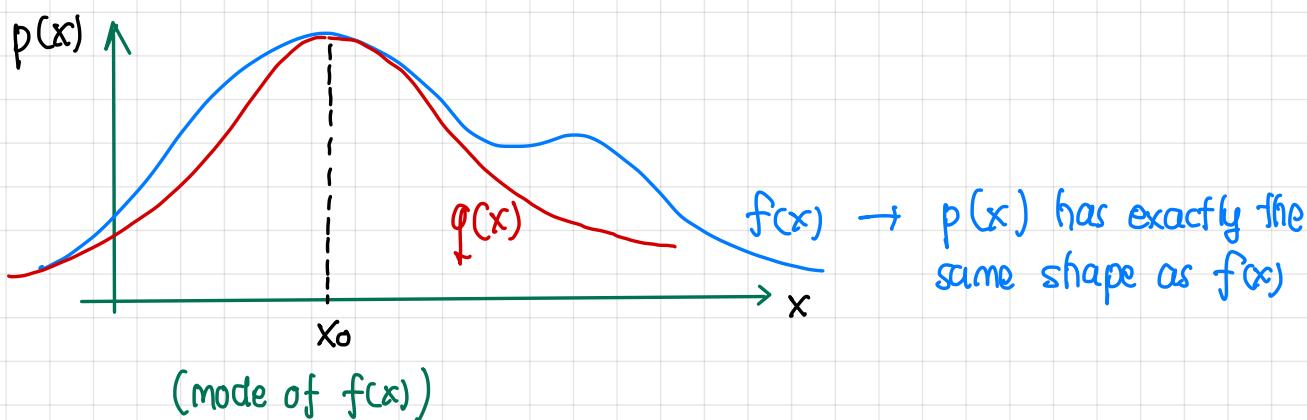
Mean of $q(x)$ = mode of $f(x)$

x_0 = mode of $f(x)$

Step 2 :

Covariance of $q(x)$ = covariance computed from $f(x)$

$$\text{thus : } \frac{d}{dx} f(x) \Big|_{x=x_0} = 0$$



Step #2 : To find the approximated covariance from $f(x)$

Let : $g(x) = \log f(x)$

Taylor expansion :

$$g(x) = g(x_0) + (x-x_0) \frac{d}{dx} g(x) \Big|_{x=x_0} + \frac{(x-x_0)^2}{2} \frac{d^2}{dx^2} g(x) \Big|_{x=x_0}$$

(Idea : to use this Taylor expansion to find the covariance from $f(x)$)

where : x_0 = the mode of $f(x)$: $\frac{d}{dx} f(x) \Big|_{x=x_0} = 0$

Implying : $\frac{d}{dx} g(x) = \frac{d}{dx} \log f(x) = \frac{1}{f(x)} \frac{d}{dx} f(x) = 0$

Hence :

$$\begin{aligned} g(x) &= g(x_0) + (x-x_0) \frac{d}{dx} g(x) \Big|_{x=x_0} + \frac{(x-x_0)^2}{2} \frac{d^2}{dx^2} g(x) \Big|_{x=x_0} \\ &= g(x_0) + \frac{(x-x_0)^2}{2} \frac{d^2}{dx^2} g(x) \Big|_{x=x_0} \end{aligned}$$

$$\log f(x) = \log f(x_0) + \frac{(x-x_0)^2}{2} \frac{d^2}{dx^2} \log f(x) \Big|_{x=x_0}$$

Let $A = - \frac{d^2}{dx^2} \log f(x) \Big|_{x=x_0}$

Then : $\log f(x) = \log f(x_0) - \frac{A}{2} (x-x_0)^2$

$$f(x) = f(x_0) \exp \left(-\frac{A}{2} (x-x_0)^2 \right)$$

Therefore : The precision from $f(x) = \frac{1}{\sigma^2} = \beta = A$

σ^2 = standard deviation

If the precision of $q(x) = A$, then:

$$q(x) = \left(\frac{A}{2\pi}\right)^{1/2} \exp\left(-\frac{A}{2}(x-x_0)^2\right) \quad ; \quad A = -\frac{d^2}{dx^2} \log f(x)$$



for vector \bar{x} : $p(\bar{x}) = \frac{1}{2} f(\bar{x})$

$$\begin{aligned} q(\bar{x}) &= \frac{|A|^{1/2}}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} (\bar{x}-\bar{x}_0)^T \underset{D \times D}{A} (\bar{x}-\bar{x}_0)\right) \\ &= \underset{D \times 1}{G}(\bar{x}; \underset{D \times 1}{\bar{x}_0}, \underset{D \times D}{A^{-1}}) \end{aligned}$$

Back to the problem we have:

$$\begin{aligned} f(\bar{\omega}) &= \prod_n \sigma_n t_n (r-\sigma_n)^{r-t_n} G(\bar{\omega}; \bar{m}_0, \bar{s}_0) \\ &\approx \underset{f}{q}(\bar{\omega}) = G(\bar{\omega}; \bar{m}_N, \bar{s}_N) \end{aligned}$$

to find



Q: What is \bar{m}_N ? What is \bar{s}_N ?

Recall from Lecture Note 18:

$$\begin{aligned} E(\bar{\omega}) &= -\log(p(E|\bar{\omega}) p(\bar{\omega})) \\ &= -\log \left[\prod_n \sigma_n t_n (r-\sigma_n)^{r-t_n} G(\bar{\omega}; \bar{m}_0, \bar{s}_0) \right] \end{aligned}$$

Hence: $E(\bar{\omega}) = -\log f(\bar{\omega})$

If $E(\bar{\omega}) = -\log f(\bar{\omega})$:

Step 1:

$$\bar{\omega}_N \text{ (the mode of } f(\bar{\omega})) = \bar{\omega}_{MAP}$$

Because $\bar{\omega}_{MAP}$ means the $\bar{\omega}$ that indicates the highest point of $E(\bar{\omega})$ or $f(\bar{\omega})$.

Step 2:

Covariance of $q(\bar{\omega}) = \mathbb{S}_N$, where $\mathbb{S}_N = \mathbb{A}^{-1}$

$$\mathbb{A} = -\nabla\nabla \log f(\bar{\omega})$$

$$= \nabla\nabla (-\log f(\bar{\omega})) = \nabla\nabla E(\bar{\omega}) \Big|_{\bar{\omega}=\bar{\omega}_{MAP}}$$

$$= \mathbb{H} = \mathbb{S}_0^{-1} + \underbrace{\Phi^T R \Phi}_{M \times M \quad M \times N \quad N \times N \quad N \times M} \quad (\text{see lecture note 18})$$

or: Textbook Eq. (4.143) page 218

Therefore:

$$f(\bar{\omega}) = p(E(\bar{\omega}) | \bar{\omega}) \approx q(\bar{\omega}) = G(\bar{\omega}; \bar{\omega}_{MAP}, \mathbb{S}_N)$$
$$= \frac{1}{(2\pi)^{D/2}} \exp \left(-\frac{1}{2} (\bar{\omega} - \bar{\omega}_{MAP})^T \mathbb{H}^{-1} (\bar{\omega} - \bar{\omega}_{MAP}) \right) \quad D \times D$$

Finally:

$$p(\bar{\omega} | \bar{\epsilon}) = \frac{p(\bar{\epsilon} | \bar{\omega}) p(\bar{\omega})}{p(\bar{\epsilon})} = q(\bar{\omega})$$

Note: $\int q(\bar{\omega}) d\bar{\omega} = 1$